

На правах рукописи

Исаченко Роман Владимирович

СНИЖЕНИЕ РАЗМЕРНОСТИ ПРОСТРАНСТВА
В ЗАДАЧАХ ДЕКОДИРОВАНИЯ СИГНАЛОВ

05.13.17 — Теоретические основы информатики

АВТОРЕФЕРАТ
диссертации на соискание ученой степени
кандидата физико-математических наук

Москва — 2021

Работа выполнена на Кафедре интеллектуальных систем Федерального государственного автономного образовательного учреждения высшего образования «Московский физико-технический институт (национальный исследовательский университет)».

Научный руководитель:

Стрижов Вадим Викторович

доктор физико-математических наук, Федеральный исследовательский центр «Информатика и управление» Российской академии наук, отдел интеллектуальных систем, ведущий научный сотрудник.

Официальные оппоненты:

Чуличков Алексей Иванович

доктор физико-математических наук, профессор, Федеральное государственное бюджетное образовательное учреждение высшего образования «Московский государственный университет имени М. В. Ломоносова», профессор кафедры математического моделирования и информатики физического факультета.

Зайцев Алексей Алексеевич

кандидат физико-математических наук, Автономная некоммерческая образовательная организация высшего образования «Сколковский институт науки и технологий», руководитель лаборатории в Центре по научным и инженерным вычислительным технологиям для задач с большими массивами данных.

Ведущая организация:

Федеральное государственное автономное образовательное учреждение высшего образования «Санкт-Петербургский национальный исследовательский университет информационных технологий, механики и оптики».

Защита состоится **6 февраля 2020 года в 13:00** на заседании диссертационного совета Д 002.073.05 при Федеральном исследовательском центре «Информатика и управление» Российской академии наук (ФИЦ ИУ РАН) по адресу: 119333, г. Москва, ул. Вавилова, д. 40.

С диссертацией можно ознакомиться в библиотеке Федерального государственного учреждения Федеральный исследовательский центр «Информатика и управление» Российской академии наук и на сайте <http://www.frccsc.ru/>

Автореферат разослан

декабря 2021 года.

И. о. ученого секретаря

диссертационного совета Д 002.073.05

д.т.н.

И. А. Матвеев

Общая характеристика работы

Актуальность темы. В работе исследуется проблема снижения размерности пространства при решении задачи декодирования сигналов. Процесс декодирования заключается в восстановлении зависимости между двумя гетерогенными наборами данных. Прогностическая модель предсказывает отклик на входной исходный сигнал.

Исходное описание данных является избыточным. При высокой мультикорреляции в признаковом пространстве финальная прогностическая модель оказывается неустойчивой. Для построения простой, устойчивой модели применяются методы снижения размерности пространства (Motrenko: 2018, Chun: 2010, Mehmood: 2012) и выбора признаков (Katrutsa: 2017, Li: 2017).

В работе решается задача декодирования с векторной целевой переменной. Пространство целевых сигналов обладает избыточной размерностью. Методы снижения размерности, не учитывающие зависимости в целевом пространстве, не являются адекватными. При предсказании векторной целевой переменной анализируется структура целевого пространства. Предложены методы, которые учитывают зависимости как в пространстве исходных объектов, так и в пространстве целевой переменной. Предлагается отобразить пространства исходных и целевых сигналов в скрытые подпространства меньшей размерности. Для построения оптимальной модели предлагаются методы согласования скрытых пространств (Wold: 1975, Rosipal: 2005, Eliseyev: 2017). Предложенные методы позволяют учесть регрессионную компоненту между исходным и целевым сигналами, а также авторегрессионную компоненту целевого сигнала.

Методы снижения размерности пространства понижают размерность исходного пространства объектов, и, как следствие, сложность модели существенно снижается (Tipping: 1999, Hotelling: 1992). Алгоритмы снижения размерности находят оптимальные комбинации исходных признаков. Если число таких комбинаций существенно меньше, чем число исходных признаков, то полученное представление снижает размерность. Цель снижения размерности — получение наиболее репрезентативных и информативных комбинаций признаков для решения задачи.

Выбор признаков является частным случаем снижения размерности пространства (Katrutsa: 2015, 2017). Найденные комбинации признаков являются подмножеством исходных признаков. Таким образом отсеиваются шумовые неинформативные признаки. Рассматриваются два типа методов выбора признаков (Li: 2017, Rodriguez-Lujan: 2010, Friedman: 2001). Первый тип методов не зависит от последующей прогностической модели. Признаки отбираются на основе свойств исходных пространств, а не на основе свойств модели. Второй тип методов отбирает признаки с учётом знания о прогностической модели.

После нахождения оптимального представления данных с помощью снижения размерности, ставится задача нахождения оптимальной метрики в скрытом пространстве объектов (Wang: 2017, Davis: 2007, Kulis: 2012, Yang: 2006,

Weinberger: 2009). В случае евклидова пространства естественным выбором метрики оказывается квадратичная норма. Задача метрического обучения заключается в нахождении оптимальной метрики, связывающей объекты.

В качестве прикладной задачи анализируется задача построения нейрокомпьютерного интерфейса (Wolpaw: 2000, Allison: 2007). Цель состоит в извлечении информации из сигналов мозговой активности (Nagel: 2018, Zhang: 2020, Chiarelli: 2018). В качестве исходных сигналов выступают сигналы электроэнцефалограммы или электрокортикограммы. Целевым сигналом является траектория движения конечности индивидуума. Необходимо построить адекватную и эффективную модель декодирования исходного сигнала в целевой сигнал. Пространство частотных характеристик мозговых сигналов и авторегрессионное пространство целевых сигналов являются чрезвычайно избыточными (Eliseyev: 2011, 2013). Построение модели без учёта имеющихся зависимостей приводит к неустойчивости модели.

В диссертации решается задача декодирования с векторной целевой переменной. Для построения оптимальной модели декодирования сигналов предлагаются методы выбора согласованных моделей с проекцией в скрытое пространство. Исходные и целевые сигналы проецируются в пространство существенно меньшей размерности. Для связи проекций исходного и целевого сигналов предлагаются методы согласования. Рассматриваются гетерогенные наборы сигналов, природа источников измерений различны. Рассматриваются как линейные методы декодирования, так и их нелинейные обобщения. Доказаны теоремы об оптимальности предложенных методов выбора моделей.

Цели работы.

1. Исследовать свойства решения задачи декодирования сигналов с векторной целевой переменной.
2. Предложить методы снижения размерности пространства, учитывающие зависимости как в пространстве исходных сигналов, так и в целевом пространстве.
3. Предложить процедуру выбора признаков для задачи декодирования сигналов.
4. Исследовать свойства линейных и нелинейных моделей для решения поставленной модели. Получить теоретические оценки оптимальности моделей.
5. Провести вычислительные эксперименты для проверки адекватности предложенных методов.

Основные положения, выносимые на защиту.

1. Исследована проблема снижения размерности сигналов в коррелированных пространствах высокой размерности. Предложены методы декодирования сигналов, учитывающие зависимости как в исходном, так и в целевом пространстве сигналов.

2. Доказаны теоремы об оптимальности предлагаемых методов декодирования сигналов. Предлагаемые методы выбирают согласованные модели в случае избыточной размерности описания данных.
3. Предложены методы выбора признаков, учитывающие зависимости как в исходном, так и в целевом пространстве. Предложенные методы доставляют устойчивые и адекватные решения в пространствах высокой размерности.
4. Предложены нелинейные методы согласования скрытых пространств для данных со сложноорганизованной целевой переменной. Предложен метод выбора активных параметров для оптимизации нелинейной модели. Исследованы свойства предлагаемого метода.
5. Предложен алгоритм метрического обучения для временных рядов с процедурой их выравнивания.
6. Предложен ряд моделей для прогнозирования гетерогенных наборов сигналов для задачи построения нейрокомпьютерных интерфейсов. Проведены вычислительные эксперименты, подтверждающие адекватность моделей.

Методы исследования. Для достижения поставленных целей используются линейные и нелинейные методы регрессионного анализа. Для анализа временных рядов используются классические авторегрессионные методы. Для извлечения признаков используются частотные характеристики временного ряда. Для построения скрытого пространства используются линейные методы снижения размерности пространства, их нелинейные модификации, а также нейросетевые методы. Для выбора признаков наряду с классическими методами, используются методы, основанные на решении задачи квадратичного программирования. Для построения метрического пространства используются методы условной выпуклой оптимизации.

Научная новизна. Предложены методы построения моделей декодирования сигналов, учитывающие структуры пространств исходных и целевых переменных. Предложены методы проекции сигналов в скрытое пространство, а также процедуры согласования образов. Предложены методы выбора признаков с помощью квадратичного программирования. Предложен метод выбора активных параметров нелинейной модели с помощью выбора признаков. Предложены методы построения оптимального метрического пространства для задачи анализа временных рядов.

Теоретическая значимость. Доказаны теоремы об оптимальности предлагаемых моделей декодирования сигналов. Доказаны теоремы о корректности рассматриваемых согласованных моделей проекций в скрытое пространство. Доказаны теоремы о достижении точки равновесия для предлагаемых методов выбора признаков.

Практическая значимость. Предложенные в работе методы предназначены для декодирования множества временных рядов сигналов электрокортикограмм, а также нестационарных временных рядов; выбора оптимальных частотных характеристик сигналов; выбора наиболее информативных параметров модели; классификации и кластеризации временных рядов физической активности.

Степень достоверности и апробация работы. Достоверность результатов подтверждена математическими доказательствами, экспериментальной проверкой результатов предлагаемых методов на реальных данных, публикациями результатов в рецензируемых научных изданиях, в том числе рекомендованных ВАК. Результаты работы докладывались и обсуждались на следующих научных конференциях.

1. Р. В. Исаченко. Метрическое обучение в задачах мультиклассовой классификации временных рядов. *Международная научная конференция «Ломоносов»*, 2016.
2. R. G. Neychev, A. P. Motrenko, R. V. Isachenko, A. S. Inyakin, and V. V. Strijov. Multimodal forecasting multiscale time series in internet of things. *Международная научная конференция «11th International Conference on Intelligent Data Processing: Theory and Applications»*, 2016.
3. Р. В. Исаченко, И. Н. Жариков, и А. М. Бочкарёв. Локальные модели для классификации объектов сложной структуры. *Всероссийская научная конференция «Математические методы распознавания образов»*, 2017.
4. R. V. Isachenko and V. V. Strijov. Dimensionality reduction for multicorrelated signal decoding with projections to latent space. *Международная научная конференция «12th International Conference on Intelligent Data Processing: Theory and Applications»*, 2018.
5. Р. В. Исаченко, В. В. Стрижов. Снижение размерности в задаче декодирования временных рядов. *Международная научная конференция «13th International Conference on Intelligent Data Processing: Theory and Applications»*, 2020.

Публикации по теме диссертации. Основные результаты по теме диссертации изложены в 6 печатных изданиях, 5 из которых изданы в журналах, рекомендованных ВАК.

Структура и объем работы. Диссертация состоит из оглавления, введения, 6 глав, заключения, списка иллюстраций, списка таблиц, списка основных обозначений и списка литературы из 110 наименований. Основной текст занимает 121 страниц.

Личный вклад. Все приведенные результаты, кроме отдельно оговоренных случаев, получены диссертантом лично при научном руководстве д.ф.-м.н. В. В. Стрижова.

Основное содержание работы

Во **введении** обоснована актуальность диссертационной работы, сформулированы цели и методы исследования, обоснована научная новизна, теоретическая и практическая значимости полученных результатов.

В **главе 1** ставится общая задача декодирования временных рядов. Ставится задача построения оптимальной линейной регрессионной модели декодирования.

Пусть $\mathbb{X} \subset \mathbb{R}^n$ — пространство исходной переменной, $\mathbb{Y} \subset \mathbb{R}^r$ — пространство целевой переменной. Пусть задано множество объектов $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^m$, где $\mathbf{x}_i \in \mathbb{X}$ — исходный объект, $\mathbf{y}_i \in \mathbb{Y}$ — целевой объект.

Обозначим за $\mathbf{X} \in \mathbb{R}^{m \times n}$ матрицу исходной переменной, за $\mathbf{Y} \in \mathbb{R}^{m \times r}$ матрицу целевой переменной:

$$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_m]^\top = [\mathbf{x}_1, \dots, \mathbf{x}_n]; \quad \mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_m]^\top = [\boldsymbol{\nu}_1, \dots, \boldsymbol{\nu}_r].$$

Столбцы $\{\mathbf{x}_j\}_{j=1}^n$ матрицы \mathbf{X} являются исходными признаками, столбцы $\{\boldsymbol{\nu}_j\}_{j=1}^r$ матрицы \mathbf{Y} являются целевыми векторами.

Предполагается, что между исходным объектом \mathbf{x} и целевым объектом \mathbf{y} существует зависимость. Требуется построить прогностическую модель $\mathbf{f} : \mathbb{X} \rightarrow \mathbb{Y}$.

Задача восстановления регрессионной зависимости состоит в нахождении оптимальной модели \mathbf{f}^* по заданным матрицам \mathbf{X} и \mathbf{Y} :

$$\mathbf{f}^* = \arg \min_{\mathbf{f}} \mathcal{L}(\mathbf{f}, \mathbf{X}, \mathbf{Y}). \quad (1)$$

Для сужения пространства поиска моделей будем рассматривать параметрические модели $\mathbf{f}(\mathbf{x}, \boldsymbol{\Theta})$, где $\boldsymbol{\Theta}$ — *параметры модели*. Таким образом между объектами \mathbf{x} и \mathbf{y} существует зависимость вида

$$\mathbf{y} = \mathbf{f}(\mathbf{x}, \boldsymbol{\Theta}) + \boldsymbol{\varepsilon},$$

где \mathbf{f} — параметрическая прогностическая модель, $\boldsymbol{\Theta}$ — параметры модели, $\boldsymbol{\varepsilon} \in \mathbb{R}^m$ — вектор регрессионных остатков.

Задача (1) сводится к задаче поиска оптимальных параметров

$$\boldsymbol{\Theta}^* = \arg \min_{\boldsymbol{\Theta}} \mathcal{L}(\boldsymbol{\Theta}, \mathbf{X}, \mathbf{Y}).$$

Предположим, что зависимость $\mathbf{f}(\mathbf{x}, \boldsymbol{\Theta})$ линейная:

$$\mathbf{y} = \mathbf{f}(\mathbf{x}, \boldsymbol{\Theta}) + \boldsymbol{\varepsilon} = \boldsymbol{\Theta}^\top \mathbf{x} + \boldsymbol{\varepsilon}, \quad (2)$$

где $\boldsymbol{\Theta} \in \mathbb{R}^{n \times r}$ — матрица параметров модели.

Оптимальные параметры $\boldsymbol{\Theta}$ определяются минимизацией функции ошибки:

$$\mathcal{L}(\boldsymbol{\Theta}, \mathbf{X}, \mathbf{Y}) = \left\| \begin{matrix} \mathbf{Y} & - & \mathbf{X} \cdot \boldsymbol{\Theta} \\ m \times r & & m \times n \quad r \times n \end{matrix} \right\|_2^2 \rightarrow \min_{\boldsymbol{\Theta}}. \quad (3)$$

Решением (3) является следующая матрица:

$$\Theta = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}.$$

Если существует вектор $\alpha \neq \mathbf{0}_n$ такой, что $\mathbf{X}\alpha = \mathbf{0}_m$, то добавление α к любому столбцу матрицы Θ не меняет значение функции потерь $\mathcal{L}(\Theta, \mathbf{X}, \mathbf{Y})$. В этом случае матрица $\mathbf{X}^T \mathbf{X}$ близка к сингулярной и не обратима.

Задача декодирования сигналов состоит в восстановлении регрессионной зависимости (1) между наборами гетерогенных сигналов.

Пусть имеется два множества временных рядов $\mathcal{S}_x = \{\mathbf{s}_x^i\}_{i=1}^m$ и $\mathcal{S}_y = \{\mathbf{s}_y^i\}_{i=1}^r$, состоящие из m и r временных рядов соответственно. Первое множество \mathcal{S}_x является множеством временных рядов m исходных сигналов. Второе множество \mathcal{S}_y является множеством временных рядов r целевых сигналов. Каждый временной ряд $\mathbf{s} = (s_1, s_2, \dots, s_T)$ является последовательностью измерений некоторой величины в течение времени.

Определение 1. *Временное представление* $\mathbf{x}_t = ([\mathbf{s}_x^1]_t, \dots, [\mathbf{s}_x^m]_t) \in \mathbb{R}^m$ состоит из измерений временных рядов исходных сигналов в момент времени t . Аналогично временное представление $\mathbf{y}_t = ([\mathbf{s}_y^1]_t, \dots, [\mathbf{s}_y^r]_t) \in \mathbb{R}^r$ состоит из измерений временных рядов целевых сигналов в момент времени t .

Определение 2. Определим *представление предыстории* длины h для момента времени t множества временных рядов исходных сигналов \mathcal{S}_x как совокупность представлений $\mathbf{X}_{t,h} = [\mathbf{x}_{t-h+1}, \dots, \mathbf{x}_t]^T \in \mathbb{R}^{h \times m}$. Аналогично определим представление предыстории длины h для момента времени t множества временных рядов целевых сигналов \mathcal{S}_y как совокупность представлений $\mathbf{Y}_{t,h} = [\mathbf{y}_{t-h+1}, \dots, \mathbf{y}_t]^T \in \mathbb{R}^{h \times r}$.

Определение 3. Определим *представление горизонта прогнозирования* длины p для момента времени t множества временных рядов исходных сигналов \mathcal{S}_x как совокупность представлений $\mathbf{X}_{t,p} = [\mathbf{x}_{t+1}, \dots, \mathbf{x}_{t+p}]^T \in \mathbb{R}^{p \times m}$. Аналогично определим представление горизонта прогнозирования длины p для момента времени t множества временных рядов целевых сигналов \mathcal{S}_y как совокупность представлений $\mathbf{Y}_{t,p} = [\mathbf{y}_{t+1}, \dots, \mathbf{y}_{t+p}]^T \in \mathbb{R}^{p \times r}$.

Определение 4. Прогностическая модель $\mathbf{f}_x^{\text{AR}} : \mathbb{R}^{h \times m} \rightarrow \mathbb{R}^{p \times m}$ является *авто-регрессионной моделью*, которая по представлению предыстории $\mathbf{X}_{t,h}$ множества временных рядов исходных сигналов \mathcal{S}_x предсказывает представление горизонта прогнозирования $\mathbf{X}_{t,p}$ множества временных рядов исходных сигналов \mathcal{S}_x . Аналогично вводится прогностическая модель $\mathbf{f}_y^{\text{AR}} : \mathbb{R}^{h \times r} \rightarrow \mathbb{R}^{p \times r}$ для множества целевых сигналов \mathcal{S}_y .

Определение 5. Определим задачу *регрессионного декодирования* как задачу построения прогностической модели $\mathbf{f}_{xy}^{\text{R}} : \mathbb{R}^{h \times m} \rightarrow \mathbb{R}^{p \times r}$, которая по представлению предыстории $\mathbf{X}_{t,h}$ множества временных рядов исходных сигналов \mathcal{S}_x предсказывает представление горизонта прогнозирования $\mathbf{Y}_{t,p}$ множества временных рядов целевых сигналов \mathcal{S}_y .

Определение 6. Общая задача декодирования состоит в построении прогностической модели $\mathbf{f}_{\mathbf{xy}} : \mathbb{R}^{h_x \times m} \times \mathbb{R}^{h_y \times r} \rightarrow \mathbb{R}^{p \times r}$, которая по представлениям предыстории \mathbf{X}_{t,h_x} и \mathbf{Y}_{t,h_y} временных рядов исходных и целевых сигналов предсказывает представление горизонта прогнозирования $\mathbf{Y}_{t,r}$ временных рядов целевых сигналов.

Задача декодирования временных рядов декомпозируется на следующие подзадачи.

- Порождение признакового пространства. Процедура порождения признакового пространства может быть основана на экспертных знаниях или же являться моделью машинного обучения.
- Снижение размерности пространства или выбор признаков. Исходные временные ряды, а также порожденное признаковое пространство оказывается избыточным, что приводит к избыточности и неустойчивости модели.
- Построение модели. После нахождения оптимального низкоразмерного представления исходных данных ставится задача выбора оптимальной модели декодирования.

Метод главных компонент для задачи декодирования. Метод главных компонент (РСА) находит низкоразмерное представление матрицы $\mathbf{X} = \mathbf{T}\mathbf{P}$, такое что новое представление $\mathbf{T} \in \mathbb{R}^{m \times l}$ содержит максимальную долю дисперсии исходной матрицы. При этом матрица отображения $\mathbf{P} \in \mathbb{R}^{l \times n}$ ($\mathbf{P}\mathbf{P}^\top = \mathbf{I}$) содержит правые собственные вектора матрицы ковариаций $\mathbf{X}^\top \mathbf{X}$.

После нахождения матрицы отображения \mathbf{P} задача (3) принимает вид

$$\mathcal{L}(\mathbf{B}, \mathbf{T}, \mathbf{Y}) = \left\| \begin{matrix} \mathbf{Y} & - & \mathbf{T} \cdot \mathbf{B} \\ m \times r & & m \times l \quad l \times r \end{matrix} \right\|_2^2 \rightarrow \min_{\mathbf{B}}.$$

Модель прогнозирования (2) в случае снижения размерности с помощью РСА принимает вид:

$$\mathbf{y} = \mathbf{B}^\top \mathbf{t} + \boldsymbol{\varepsilon} = \mathbf{B}^\top \mathbf{P} \mathbf{x} + \boldsymbol{\varepsilon} = \boldsymbol{\Theta} \mathbf{x} + \boldsymbol{\varepsilon}, \text{ где } \boldsymbol{\Theta} = \mathbf{B}^\top \mathbf{P}.$$

Метод частичных наименьших квадратов для задачи декодирования. Основным недостатком метода РСА является отсутствие учёта взаимосвязи между исходными признаками χ_j и целевыми векторами ν_j . Метод частичных наименьших квадратов (PLS) проецирует матрицу исходных объектов \mathbf{X} и матрицу целевых объектов \mathbf{Y} в скрытое пространство малой размерностью l ($l < n$). Метод PLS находит в скрытом пространстве матрицы $\mathbf{T}, \mathbf{U} \in \mathbb{R}^{m \times l}$, которые лучше всего описывают исходные матрицы \mathbf{X} и \mathbf{Y} . При этом PLS максимизирует ковариацию между столбцами матриц \mathbf{T} и \mathbf{U} соответственно. Принцип работы

метода показан на следующей коммутативной диаграмме:

$$\begin{array}{ccc}
 \mathbf{x} \in \mathbb{R}^n & \xrightarrow{\mathbf{f}} & \mathbf{y} \in \mathbb{R}^r \\
 & \swarrow \mathbf{P} \quad \searrow \mathbf{Q} & \\
 & \mathbf{t}, \mathbf{u} \in \mathbb{R}^\ell & \\
 & \nwarrow \mathbf{W} \quad \nearrow \mathbf{C} &
 \end{array}$$

Матрица исходных объектов \mathbf{X} и матрица целевых объектов \mathbf{Y} проецируются на скрытое пространство следующим образом:

$$\begin{aligned}
 \mathbf{X}_{m \times n} &= \mathbf{T}_{m \times l} \cdot \mathbf{P}_{l \times n} + \mathbf{E}_x_{m \times n} = \sum_{k=1}^l \boldsymbol{\tau}_k \cdot \mathbf{p}_k^\top + \mathbf{E}_x_{m \times n}, \\
 \mathbf{Y}_{m \times r} &= \mathbf{U}_{m \times l} \cdot \mathbf{Q}_{l \times r} + \mathbf{E}_y_{m \times r} = \sum_{k=1}^l \boldsymbol{\nu}_k \cdot \mathbf{q}_k^\top + \mathbf{E}_y_{m \times r}.
 \end{aligned}$$

Здесь \mathbf{T} и \mathbf{U} — образы исходных матриц в скрытом пространстве, причём столбцы матрицы \mathbf{T} ортогональны; \mathbf{P} и \mathbf{Q} — матрицы перехода; \mathbf{E}_x и \mathbf{E}_y — матрицы остатков. Метод PLS восстанавливает линейную зависимость между столбцами матриц \mathbf{T} и \mathbf{U}

$$\mathbf{U} \approx \mathbf{T}\mathbf{B}, \quad \mathbf{B} = \text{diag}(\beta_k), \quad \beta_k = \boldsymbol{\nu}_k^\top \boldsymbol{\tau}_k / (\boldsymbol{\tau}_k^\top \boldsymbol{\tau}_k),$$

где $\{\boldsymbol{\tau}_k\}_{k=1}^l$, $\{\boldsymbol{\nu}_k\}_{k=1}^l$ — столбцы матриц \mathbf{T} и \mathbf{U} соответственно.

Метод решает следующую оптимизационную задачу:

$$\max_{\|\mathbf{p}\|_2=\|\mathbf{q}\|_2=1} [\text{cov}(\mathbf{X}\mathbf{p}, \mathbf{Y}\mathbf{q})^2] = \max_{\mathbf{p}, \mathbf{q}} \frac{\mathbf{p}^\top \mathbf{X}^\top \mathbf{Y} \mathbf{q}}{\sqrt{\mathbf{p}^\top \mathbf{p}} \sqrt{\mathbf{q}^\top \mathbf{q}}}. \quad (4)$$

Канонический анализ корреляций для задачи декодирования. Оптимизационная задача канонического корреляционного анализа (ССА) похожа на оптимизационную задачу PLS (4) с той лишь разницей, что вместо максимизации ковариации ССА максимизирует корреляцию:

$$\max_{\|\mathbf{p}\|_2=\|\mathbf{q}\|_2=1} [\text{corr}(\mathbf{X}\mathbf{p}, \mathbf{Y}\mathbf{q})^2] = \max_{\mathbf{p}, \mathbf{q}} \frac{\mathbf{p}^\top \mathbf{X}^\top \mathbf{Y} \mathbf{q}}{\sqrt{\mathbf{p}^\top \mathbf{X}^\top \mathbf{X} \mathbf{p}} \sqrt{\mathbf{q}^\top \mathbf{Y}^\top \mathbf{Y} \mathbf{q}}}.$$

Метод ДеерССА максимизирует корреляцию между представлениями, полученными на выходе нейросети:

$$\begin{aligned}
 &\max_{\|\mathbf{p}\|_2=\|\mathbf{q}\|_2=1} [\text{corr}(\mathbf{G}_x(\mathbf{X}, \mathbf{W}_x) \cdot \mathbf{p}, \mathbf{G}_y(\mathbf{Y}, \mathbf{W}_y) \cdot \mathbf{q})^2] = \\
 &= \max_{\mathbf{p}, \mathbf{q}} \frac{\mathbf{p}^\top \mathbf{G}_x(\mathbf{X}, \mathbf{W}_x)^\top \mathbf{G}_y(\mathbf{Y}, \mathbf{W}_y) \mathbf{q}}{\sqrt{\mathbf{p}^\top \mathbf{G}_x(\mathbf{X}, \mathbf{W}_x)^\top \mathbf{G}_x(\mathbf{X}, \mathbf{W}_x) \mathbf{p}} \sqrt{\mathbf{q}^\top \mathbf{G}_y(\mathbf{Y}, \mathbf{W}_y)^\top \mathbf{G}_y(\mathbf{Y}, \mathbf{W}_y) \mathbf{q}}}.
 \end{aligned}$$

Здесь $\mathbf{G}_x(\mathbf{X}, \mathbf{W}_x)$ и $\mathbf{G}_y(\mathbf{Y}, \mathbf{W}_y)$ — нелинейные проекции исходных и целевых объектов. С использованием нейросетевых функций модель декодирования способна учитывать существенно нелинейные зависимости как в исходном пространстве, так и в целевом пространстве.

В **главе 2** приводится формальная постановка задачи построения согласованных моделей декодирования. Вводятся понятия скрытого пространства и процедуры согласования образов. Доказываются теоремы о выборе оптимальной модели декодирования.

Предположение 1. Пусть пространства \mathbb{X} и \mathbb{Y} имеют избыточную размерность. В простейшем случае такими многообразия могут являться вложениями или линейными подпространствами.

Определение 7. Назовём пространство $\mathbb{T} \subset \mathbb{R}^l$ *скрытым пространством* для пространства $\mathbb{X} \subset \mathbb{R}^n$ ($l \leq n$), если существуют *функция кодирования* $\varphi_x : \mathbb{X} \rightarrow \mathbb{T}$ и *функция декодирования* $\psi_x : \mathbb{T} \rightarrow \mathbb{X}$, такие что

$$\text{для любого } \mathbf{x} \in \mathbb{X} \quad \text{существует } \mathbf{t} \in \mathbb{T} : \psi_x(\varphi_x(\mathbf{x})) = \psi_x(\mathbf{t}) = \mathbf{x}.$$

Аналогично введём определение *скрытого пространства* $\mathbb{U} \subset \mathbb{R}^s$ для целевого пространства \mathbb{Y} , *функции кодирования* $\varphi_y : \mathbb{Y} \rightarrow \mathbb{U}$ и *декодирования* $\psi_y : \mathbb{U} \rightarrow \mathbb{Y}$, такие что

$$\text{для любого } \mathbf{y} \in \mathbb{Y} \quad \text{существует } \mathbf{u} \in \mathbb{U} : \psi_y(\varphi_y(\mathbf{y})) = \psi_y(\mathbf{u}) = \mathbf{y}.$$

Образы матрицы исходных объектов \mathbf{X} и матрицы целевых объектов \mathbf{Y} в скрытых пространствах \mathbb{T} и \mathbb{U} имеют вид

$$\begin{aligned} \mathbf{T} &= \varphi_x(\mathbf{X}) = [\mathbf{t}_1, \dots, \mathbf{t}_n]^\top = [\boldsymbol{\tau}_1, \dots, \boldsymbol{\tau}_l], \\ \mathbf{U} &= \varphi_y(\mathbf{Y}) = [\mathbf{u}_1, \dots, \mathbf{u}_s]^\top = [\boldsymbol{\nu}_1, \dots, \boldsymbol{\nu}_s]. \end{aligned}$$

Определение 8. Будем говорить, что скрытые пространства \mathbb{T} и \mathbb{U} являются *согласованными*, если существует *функция связи* $\mathbf{h} : \mathbb{T} \rightarrow \mathbb{U}$, такая что

$$\text{для любого } \mathbf{u} \in \mathbb{U} \quad \text{существует } \mathbf{t} \in \mathbb{T} : \mathbf{u} = \mathbf{h}(\mathbf{t}).$$

Предположение 2. Предположим, что в задаче прогнозирования (1) пространства \mathbb{T} и \mathbb{U} являются скрытыми для пространств \mathbb{X} и \mathbb{Y} соответственно. Предположим также, что для данных скрытых пространств \mathbb{T} и \mathbb{U} существует функция связи $\mathbf{h} : \mathbb{T} \rightarrow \mathbb{U}$. Тогда выполнено

$$\text{для любого } \mathbf{y} \in \mathbb{Y} \quad \text{существует } \mathbf{x} \in \mathbb{X} : \mathbf{y} = \psi_y(\mathbf{u}) = \psi_y(\mathbf{h}(\mathbf{t})) = \psi_y(\mathbf{h}(\varphi_x(\mathbf{x}))),$$

и общая схема задачи поиска согласованной модели декодирования принимает

вид следующей коммутативной диаграммы:

$$\begin{array}{ccc}
 \mathbb{X} \subset \mathbb{R}^n & \xrightarrow{\mathbf{f}} & \mathbb{Y} \subset \mathbb{R}^r \\
 \varphi_{\mathbf{x}} \left(\begin{array}{c} \uparrow \\ \downarrow \end{array} \right) \psi_{\mathbf{x}} & & \psi_{\mathbf{y}} \left(\begin{array}{c} \uparrow \\ \downarrow \end{array} \right) \varphi_{\mathbf{y}} \\
 \mathbb{T} \subset \mathbb{R}^l & \xrightarrow{\mathbf{h}} & \mathbb{U} \subset \mathbb{R}^s
 \end{array} \tag{5}$$

Определение 9. Согласно диаграмме (5), определим *согласованную* модель декодирования $\mathbf{f} : \mathbb{X} \rightarrow \mathbb{Y}$ как суперпозицию

$$\mathbf{f} = \psi_{\mathbf{y}} \circ \mathbf{h} \circ \varphi_{\mathbf{x}}.$$

Для поиска оптимальных параметров функций кодирования $\varphi_{\mathbf{x}}$ и $\varphi_{\mathbf{y}}$, декодирования $\psi_{\mathbf{x}}$ и $\psi_{\mathbf{y}}$, а также функции связи \mathbf{h} ставится задача максимизации функции согласования скрытых векторов

$$g : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}, \quad g(\boldsymbol{\tau}, \boldsymbol{\nu}) \rightarrow \max_{\varphi_{\mathbf{x}}, \varphi_{\mathbf{y}}, \mathbf{h}}.$$

Каждая пара скрытых векторов $\boldsymbol{\tau}, \boldsymbol{\nu}$ ищется последовательно.

Метод главных компонент. Функции кодирования $\varphi_{\mathbf{x}} : \mathbb{R}^n \rightarrow \mathbb{R}^l$ и декодирования $\psi_{\mathbf{x}} : \mathbb{R}^l \rightarrow \mathbb{R}^n$ имеют вид

$$\varphi_{\mathbf{x}}(\mathbf{X}) = \underset{m \times n}{\mathbf{X}} \cdot \underset{n \times l}{\mathbf{P}}^{\top}, \quad \psi_{\mathbf{x}}(\mathbf{T}) = \underset{m \times l}{\mathbf{T}} \cdot \underset{l \times n}{\mathbf{P}},$$

где $\mathbf{P} = [\mathbf{p}_1, \dots, \mathbf{p}_l]^{\top}$, при этом $\mathbf{P}\mathbf{P}^{\top} = \mathbf{I}$.

Скрытые вектора $\boldsymbol{\tau}$ строятся так, чтобы выборочная дисперсия столбцов проекций матрицы \mathbf{X} была максимальной:

$$\mathbf{p} = \arg \max_{\|\mathbf{p}\|_2=1} g(\boldsymbol{\tau}) = \arg \max_{\|\mathbf{p}\|_2=1} [\text{var}(\boldsymbol{\tau})] = \arg \max_{\|\mathbf{p}\|_2=1} [\text{var}(\mathbf{X}\mathbf{p})].$$

Метод PCA не согласует исходные переменные и целевые переменные. А именно метод PCA не находит функции кодирования $\varphi_{\mathbf{y}}$ и декодирования $\psi_{\mathbf{y}}$, а также функцию связи \mathbf{h} . При этом функция согласования скрытых векторов $g(\boldsymbol{\tau})$ зависит только от одного аргумента.

Метод частичных наименьших квадратов и канонический анализ корреляций. В методах PLS и CCA функции кодирования и декодирования имеют вид

$$\begin{aligned}
 \varphi_{\mathbf{x}}(\mathbf{X}) &= \mathbf{X}\mathbf{W}, & \varphi_{\mathbf{y}}(\mathbf{Y}) &= \mathbf{Y}\mathbf{C}, \\
 \psi_{\mathbf{x}}(\mathbf{T}) &= \mathbf{T}\mathbf{P}^{\top}, & \psi_{\mathbf{y}}(\mathbf{U}) &= \mathbf{U}\mathbf{Q}^{\top}.
 \end{aligned}$$

Функция связи \mathbf{h} имеет вид линейной модели, связывающей образы проекций в скрытом пространстве $\mathbf{u} = \mathbf{h}(\mathbf{t}) = \mathbf{B}^\top \mathbf{t}$. В данном случае схема декодирования (5) принимает вид следующей коммутативной диаграммы.

$$\begin{array}{ccc} \mathbf{X} \subset \mathbb{R}^n & \xrightarrow{\mathbf{f}} & \mathbf{Y} \subset \mathbb{R}^r \\ \mathbf{W} \left(\begin{array}{c} \uparrow \\ \downarrow \end{array} \right) \mathbf{P} & & \mathbf{Q} \left(\begin{array}{c} \uparrow \\ \downarrow \end{array} \right) \mathbf{C} \\ \mathbf{T} \subset \mathbb{R}^\ell & \xrightarrow{\mathbf{B}} & \mathbf{U} \subset \mathbb{R}^s \end{array}$$

Для метода PLS функция согласования скрытых векторов имеет вид $g(\boldsymbol{\tau}, \boldsymbol{\nu}) = \text{cov}(\boldsymbol{\tau}, \boldsymbol{\nu})$, а для метода ССА: $g(\boldsymbol{\tau}, \boldsymbol{\nu}) = \text{corr}(\boldsymbol{\tau}, \boldsymbol{\nu})$.

Нелинейный канонический анализ корреляций. Функции кодирования и декодирования являются нелинейными нейросетями вида

$$\begin{aligned} \mathbf{T} &= \boldsymbol{\varphi}_{\mathbf{x}}(\mathbf{X}) = \mathbf{W}_{\mathbf{x}}^L \sigma(\dots \sigma(\mathbf{W}_{\mathbf{x}}^2 \sigma(\mathbf{X} \mathbf{W}_{\mathbf{x}}^1)) \dots), \\ \mathbf{U} &= \boldsymbol{\varphi}_{\mathbf{y}}(\mathbf{Y}) = \mathbf{W}_{\mathbf{y}}^L \sigma(\dots \sigma(\mathbf{W}_{\mathbf{y}}^2 \sigma(\mathbf{Y} \mathbf{W}_{\mathbf{y}}^1)) \dots), \\ \mathbf{X} &= \boldsymbol{\psi}_{\mathbf{x}}(\mathbf{T}) = \mathbf{W}_{\mathbf{t}}^L \sigma(\dots \sigma(\mathbf{W}_{\mathbf{t}}^2 \sigma(\mathbf{T} \mathbf{W}_{\mathbf{t}}^1)) \dots), \\ \mathbf{Y} &= \boldsymbol{\psi}_{\mathbf{y}}(\mathbf{U}) = \mathbf{W}_{\mathbf{u}}^L \sigma(\dots \sigma(\mathbf{W}_{\mathbf{u}}^2 \sigma(\mathbf{U} \mathbf{W}_{\mathbf{u}}^1)) \dots). \end{aligned}$$

Каждая нейросеть является суперпозицией последовательных умножений на матрицы параметров и применения поэлементных функций активаций.

Требуется найти такие параметры, при которых функция согласования g достигает своего максимума:

$$g(\boldsymbol{\tau}, \boldsymbol{\nu}) \rightarrow \max_{\mathbf{W}},$$

где $\mathbf{W} = \{\mathbf{W}_{\mathbf{x}}^i, \mathbf{W}_{\mathbf{y}}^i, \mathbf{W}_{\mathbf{t}}^i, \mathbf{W}_{\mathbf{u}}^i\}_{i=1}^L$.

При использовании в качестве функции согласования корреляции между проекциями $g(\boldsymbol{\tau}, \boldsymbol{\nu}) = \text{corr}(\boldsymbol{\tau}, \boldsymbol{\nu})$ частная производная функции согласования по первому аргументу принимает вид

$$\frac{\partial \mathbf{g}(\mathbf{T}, \mathbf{U})}{\partial \mathbf{T}} = \frac{1}{m-1} \left(\boldsymbol{\Sigma}_1^{-1/2} \mathbf{U} \mathbf{V}^\top \boldsymbol{\Sigma}_2^{-1/2} \mathbf{U} - \boldsymbol{\Sigma}_1^{-1/2} \mathbf{U} \mathbf{D} \mathbf{V}^\top \boldsymbol{\Sigma}_1^{-1/2} \right),$$

где $\mathbf{U}, \mathbf{D}, \mathbf{V} = \text{SVD}(\boldsymbol{\Sigma})$, $\boldsymbol{\Sigma} = \boldsymbol{\Sigma}_1^{-1/2} \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_2^{-1/2}$, $\boldsymbol{\Sigma}_1 = \frac{1}{m-1} \mathbf{T} \mathbf{T}^\top$, $\boldsymbol{\Sigma}_2 = \frac{1}{m-1} \mathbf{U} \mathbf{U}^\top$, $\boldsymbol{\Sigma}_{12} = \frac{1}{m-1} \mathbf{T} \mathbf{U}^\top$. Аналогичное выражение имеет частная производная по второму аргументу. Полученное выражение для градиента позволяет построить эффективный алгоритм для решения задачи с использованием градиентных методов оптимизации.

Алгоритм PLS итеративно на каждом из l шагов вычисляет по одному столбцу $\boldsymbol{\tau}_k$, $\boldsymbol{\nu}_k$ матриц \mathbf{T} , \mathbf{U} и по одной строке \mathbf{p}_k , \mathbf{q}_k матриц \mathbf{P} , \mathbf{Q} соответственно.

Утверждение 1. Максимизация ковариации между векторами $\boldsymbol{\tau}_k$ и $\boldsymbol{\nu}_k$ сохраняет дисперсию столбцов матриц \mathbf{X} и \mathbf{Y} и учитывает их линейную зависимость.

Утверждение 2. Вычисленные вектора \mathbf{w}_k и \mathbf{c}_k с помощью итеративной процедуры обновления:

$$\boldsymbol{\tau}_k := \frac{\mathbf{X}_k \mathbf{w}_k}{\|\mathbf{w}_k\|}, \quad \mathbf{w}_k := \frac{\mathbf{X}_k^\top \boldsymbol{\nu}_{k-1}}{\boldsymbol{\nu}_{k-1}^\top \boldsymbol{\nu}_{k-1}}; \quad (6)$$

$$\boldsymbol{\nu}_k := \frac{\mathbf{Y}_k \mathbf{c}_k}{\|\mathbf{c}_k\|}, \quad \mathbf{c}_k := \frac{\mathbf{Y}_k^\top \boldsymbol{\tau}_k}{\boldsymbol{\tau}_k^\top \boldsymbol{\tau}_k}. \quad (7)$$

будут собственными векторами матриц $\mathbf{X}_k^\top \mathbf{Y}_k \mathbf{Y}_k^\top \mathbf{X}_k$ и $\mathbf{Y}_k^\top \mathbf{X}_k \mathbf{X}_k^\top \mathbf{Y}_k$, соответствующими максимальным собственным значениям.

$$\mathbf{w}_k \propto \mathbf{X}_k^\top \boldsymbol{\nu}_{k-1} \propto \mathbf{X}_k^\top \mathbf{Y}_k \mathbf{c}_{k-1} \propto \mathbf{X}_k^\top \mathbf{Y}_k \mathbf{Y}_k^\top \boldsymbol{\tau}_{k-1} \propto \mathbf{X}_k^\top \mathbf{Y}_k \mathbf{Y}_k^\top \mathbf{X}_k \mathbf{w}_{k-1},$$

$$\mathbf{c}_k \propto \mathbf{Y}_k^\top \boldsymbol{\tau}_k \propto \mathbf{Y}_k^\top \mathbf{X}_k \mathbf{w}_k \propto \mathbf{Y}_k^\top \mathbf{X}_k \mathbf{X}_k^\top \boldsymbol{\nu}_{k-1} \propto \mathbf{Y}_k^\top \mathbf{X}_k \mathbf{X}_k^\top \mathbf{Y}_k \mathbf{c}_{k-1},$$

где символ \propto означает равенство с точностью до мультипликативной константы.

Утверждение 3. Вычисленные вектора $\boldsymbol{\tau}_k$ и $\boldsymbol{\nu}_k$ с помощью итеративной процедуры обновления (6), (7) обладают максимальной ковариацией $\text{cov}(\boldsymbol{\tau}, \boldsymbol{\nu})$.

Для перехода на следующий шаг необходимо вычесть из матриц \mathbf{X}_k и \mathbf{Y}_k одноранговые аппроксимации $\boldsymbol{\tau}_k \mathbf{p}_k^\top$ и $\boldsymbol{\tau}_k \mathbf{q}_k^\top$

$$\begin{aligned} \mathbf{X}_{k+1} &= \mathbf{X}_k - \boldsymbol{\tau}_k \mathbf{p}_k^\top = \mathbf{X} - \sum_k \boldsymbol{\tau}_k \mathbf{p}_k^\top, \\ \mathbf{Y}_{k+1} &= \mathbf{Y}_k - \boldsymbol{\tau}_k \mathbf{q}_k^\top = \mathbf{Y} - \sum_k \boldsymbol{\tau}_k \mathbf{q}_k^\top. \end{aligned}$$

Теорема 1. В случае линейных функций декодирования $\boldsymbol{\psi}_x(\mathbf{T}) = \mathbf{T}\mathbf{P}$, $\boldsymbol{\psi}_y(\mathbf{U}) = \mathbf{U}\mathbf{Q}$ и функции согласования $g(\boldsymbol{\tau}, \boldsymbol{\nu}) = \text{cov}(\boldsymbol{\tau}, \boldsymbol{\nu})$ параметры

$$\boldsymbol{\Theta} = \mathbf{W}(\mathbf{P}^\top \mathbf{W})^{-1} \mathbf{B}\mathbf{Q}$$

являются оптимальными для модели (2).

Финальная модель является линейной, низкоразмерной в скрытом пространстве. Это снижает избыточность данных и повышает стабильность модели.

Пусть $\mathbf{f}_1(\mathbf{x}_1, \boldsymbol{\Theta}_1)$, $\mathbf{f}_2(\mathbf{x}_2, \boldsymbol{\Theta}_2)$ — линейные модели декодирования сигналов.

Определение 10. Назовём *аддитивной суперпозицией* моделей декодирования модель (2) вида

$$\mathbf{Y} = \mathbf{f}(\mathbf{x}, \boldsymbol{\Theta}) + \boldsymbol{\varepsilon} = \mathbf{f}_1(\mathbf{x}_1, \boldsymbol{\Theta}_1) + \mathbf{f}_2(\mathbf{x}_2, \boldsymbol{\Theta}_2) + \boldsymbol{\varepsilon} = \boldsymbol{\Theta}_1^\top \mathbf{x}_1 + \boldsymbol{\Theta}_2^\top \mathbf{x}_2 + \boldsymbol{\varepsilon}, \quad (8)$$

где объект $\mathbf{x} = [\mathbf{x}_1, \mathbf{x}_2] \in \mathbb{R}^n$ состоит из двух подвекторов $\mathbf{x}_1 \in \mathbb{R}^k$, $\mathbf{x}_2 \in \mathbb{R}^{n-k}$.

Утверждение 4. Оптимальная матрица параметров Θ для модели (8), доставляющая минимум функции ошибки (3), имеет вид:

$$\begin{aligned}\Theta_1 &= (\mathbf{X}_1^\top \mathbf{M}_{\mathbf{X}_2} \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \mathbf{M}_{\mathbf{X}_2} \mathbf{Y}, \\ \Theta_2 &= (\mathbf{X}_2^\top \mathbf{M}_{\mathbf{X}_1} \mathbf{X}_2)^{-1} \mathbf{X}_2^\top \mathbf{M}_{\mathbf{X}_1} \mathbf{Y},\end{aligned}$$

где $\mathbf{M}_{\mathbf{X}_1} = (\mathbf{I} - \mathbf{P}_{\mathbf{X}_1})$, $\mathbf{M}_{\mathbf{X}_2} = (\mathbf{I} - \mathbf{P}_{\mathbf{X}_2})$, $\mathbf{P}_{\mathbf{X}_1} = \mathbf{X}_1(\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top$, $\mathbf{P}_{\mathbf{X}_2} = \mathbf{X}_2(\mathbf{X}_2^\top \mathbf{X}_2)^{-1} \mathbf{X}_2^\top$.

Утверждение 5. Оптимальная подматрица Θ_2 в модели (8) является решением задачи регрессии

$$\|\mathbf{Y}_1 - \mathbf{X}_{21} \Theta_2\|^2 \rightarrow \min_{\Theta_2},$$

где $\mathbf{Y}_1 = \mathbf{M}_{\mathbf{X}_1} \mathbf{Y}$, $\mathbf{X}_{21} = \mathbf{M}_{\mathbf{X}_1} \mathbf{X}_2$.

Утверждение 6. Если в задаче (8) $\text{span}(\mathbf{X}_1) \cap \text{span}(\mathbf{X}_2) = \emptyset$, то есть столбцы матрицы \mathbf{X}_1 ортогональны столбцам матрицы \mathbf{X}_2 , то Θ_2 является решением задачи регрессии

$$\|\mathbf{Y} - \mathbf{X}_2 \Theta_2\|^2 \rightarrow \min_{\Theta_2}.$$

Утверждение 7. Ошибка аддитивной суперпозиции моделей не превышает ошибки каждой из отдельных моделей

$$\begin{aligned}\mathcal{L}(\Theta^*, \mathbf{X}, \mathbf{Y}) &\leq \mathcal{L}(\Theta_1, \mathbf{X}_1, \mathbf{Y}), \\ \mathcal{L}(\Theta^*, \mathbf{X}, \mathbf{Y}) &\leq \mathcal{L}(\Theta_2, \mathbf{X}_2, \mathbf{Y}).\end{aligned}$$

Утверждение 8. Пусть для аддитивной суперпозиции моделей (8) выполнены следующие условия

$$\mathbf{Y} \neq \mathbf{P}_{\mathbf{X}_2} \mathbf{Y}, \quad \mathbf{X}_1 \neq \mathbf{P}_{\mathbf{X}_2} \mathbf{X}_1, \quad \mathbf{Y}^\top \mathbf{M}_{\mathbf{X}_2} \mathbf{X}_1 \neq 0.$$

Тогда выполнено строгое неравенство

$$\mathcal{L}(\Theta^*, \mathbf{X}, \mathbf{Y}) < \mathcal{L}(\Theta_2, \mathbf{X}_2, \mathbf{Y}).$$

Глава 3 посвящена задаче выбора признаков для задачи декодирования сигналов. Задача выбора признаков заключается в поиске оптимального подмножества $\mathcal{A} \subset \{1, \dots, n\}$ индексов признаков среди всех возможных $2^n - 1$ вариантов. Существует взаимоднозначное отображение между подмножеством \mathcal{A} и булевым вектором $\mathbf{a} \in \{0, 1\}^n$, компоненты которого указывают, выбран ли признак. Для нахождения оптимального вектора \mathbf{a} введем функцию ошибки выбора признаков $S(\mathbf{a}, \mathbf{X}, \mathbf{Y})$:

$$\mathbf{a} = \arg \min_{\mathbf{a}' \in \{0, 1\}^n} S(\mathbf{a}', \mathbf{X}, \mathbf{Y}). \quad (9)$$

Целью выбора признаков является построение функции $S(\mathbf{a}, \mathbf{X}, \mathbf{Y})$.

Для решения задачи (9) применяется релаксация задачи (9) к непрерывной области определения $[0, 1]^n$:

$$\mathbf{z} = \arg \min_{\mathbf{z}' \in [0, 1]^n} S(\mathbf{z}', \mathbf{X}, \mathbf{Y}).$$

Решение (9) восстанавливается с помощью отсечения по порогу:

$$\mathbf{a} = [a_j]_{j=1}^n, \quad a_j = \begin{cases} 1, & z_j > \tau; \\ 0, & \text{в противном случае.} \end{cases}$$

Как только решение \mathbf{a} задачи (9) получено, задача (3) принимает вид:

$$\mathcal{L}(\Theta_{\mathcal{A}}, \mathbf{X}_{\mathcal{A}}, \mathbf{Y}) = \|\mathbf{Y} - \mathbf{X}_{\mathcal{A}}\Theta_{\mathcal{A}}\|_2^2 \rightarrow \min_{\Theta_{\mathcal{A}}}.$$

Если между столбцами матрицы исходных объектов \mathbf{X} существует линейная зависимость, то решение задачи линейной регрессии

$$\|\mathbf{v} - \mathbf{X}\boldsymbol{\theta}\|_2^2 \rightarrow \min_{\boldsymbol{\theta} \in \mathbb{R}^n}.$$

оказывается неустойчивым. Методы выбора признаков находят подмножество $\mathcal{A} \in \{1, \dots, n\}$ оптимальных столбцов матрицы \mathbf{X} .

Метод QPFS выбирает некоррелированные признаки, релевантные целевому вектору \mathbf{v} . Введем две функции: $\text{Sim}(\mathbf{X})$ контролирует избыточность между признаками, $\text{Rel}(\mathbf{X}, \mathbf{v})$ содержит релевантности между каждым признаком и целевым вектором.

Метод QPFS минимизирует следующую функцию ошибки

$$\underbrace{\mathbf{z}^T \mathbf{Q} \mathbf{z}}_{\text{Sim}} - \alpha \cdot \underbrace{\mathbf{b}^T \mathbf{z}}_{\text{Rel}} \rightarrow \min_{\substack{\mathbf{z} \in \mathbb{R}_+^n \\ \|\mathbf{z}\|_1 = 1}}. \quad (10)$$

Элементы матрицы парных взаимодействий $\mathbf{Q} \in \mathbb{R}^{n \times n}$ содержат коэффициенты попарного сходства между признаками. Вектор релевантностей признаков $\mathbf{b} \in \mathbb{R}^n$ выражает сходство между каждым признаком и целевым вектором \mathbf{v} . Нормированный вектор \mathbf{z} отражает значимость каждого признака. Параметр α позволяет контролировать компромисс между Sim и Rel:

$$\alpha = \frac{\overline{\mathbf{Q}}}{\overline{\mathbf{Q}} + \overline{\mathbf{b}}}, \quad \text{где } \overline{\mathbf{Q}} = \text{mean}(\mathbf{Q}), \quad \overline{\mathbf{b}} = \text{mean}(\mathbf{b}).$$

Для измерения сходства используется выборочный коэффициент корреляции Пирсона между парами признаков для функции Sim, и между признаками и целевым вектором для функции Rel:

$$\mathbf{Q} = [|\text{corr}(\mathbf{x}_i, \mathbf{x}_j)|]_{i,j=1}^n, \quad \mathbf{b} = [|\text{corr}(\mathbf{x}_i, \mathbf{v})|]_{i=1}^n.$$

Здесь

$$\text{corr}(\mathbf{x}, \mathbf{v}) = \frac{\sum_{i=1}^m (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{v}_i - \bar{\mathbf{v}})}{\sqrt{\sum_{i=1}^m (\mathbf{x}_i - \bar{\mathbf{x}})^2 \sum_{i=1}^m (\mathbf{v}_i - \bar{\mathbf{v}})^2}}.$$

В случае векторной целевой переменной компоненты целевой переменной могут коррелировать между собой. Предлагаются методы, учитывающие зависимости как в исходном, так и в целевом пространствах.

Агрегация релевантностей целевых векторов. Чтобы применить метод QPFS к векторному случаю ($r > 1$), релевантности признаков агрегируются по всем r компонентам целевой переменной. Вектор релевантностей \mathbf{b} агрегируется по всем компонентам целевой переменной и определяется как

$$\mathbf{b} = \left[\sum_{k=1}^r |\text{corr}(\mathbf{x}_i, \mathbf{v}_k)| \right]_{i=1}^n.$$

Недостатком такого подхода является отсутствие учёта зависимостей в столбцах матрицы \mathbf{Y} .

Симметричный учёт значимости признаков и целевых переменных. Добавим член $\text{Sim}(\mathbf{Y})$ и изменим член $\text{Rel}(\mathbf{X}, \mathbf{Y})$ следующим образом:

$$\alpha_1 \cdot \underbrace{\mathbf{z}_x^\top \mathbf{Q}_x \mathbf{z}_x}_{\text{Sim}(\mathbf{X})} - \alpha_2 \cdot \underbrace{\mathbf{z}_x^\top \mathbf{B} \mathbf{z}_y}_{\text{Rel}(\mathbf{X}, \mathbf{Y})} + \alpha_3 \cdot \underbrace{\mathbf{z}_y^\top \mathbf{Q}_y \mathbf{z}_y}_{\text{Sim}(\mathbf{Y})} \rightarrow \min_{\substack{\mathbf{z}_x \geq \mathbf{0}_n, \mathbf{1}_n^\top \mathbf{z}_x = 1 \\ \mathbf{z}_y \geq \mathbf{0}_r, \mathbf{1}_r^\top \mathbf{z}_y = 1}}. \quad (11)$$

Определим элементы матриц $\mathbf{Q}_x \in \mathbb{R}^{n \times n}$, $\mathbf{Q}_y \in \mathbb{R}^{r \times r}$ и $\mathbf{B} \in \mathbb{R}^{n \times r}$ следующим образом:

$$\mathbf{Q}_x = [|\text{corr}(\mathbf{x}_i, \mathbf{x}_j)|]_{i,j=1}^n, \quad \mathbf{Q}_y = [|\text{corr}(\mathbf{v}_i, \mathbf{v}_j)|]_{i,j=1}^r, \quad \mathbf{B} = [|\text{corr}(\mathbf{x}_i, \mathbf{v}_j)|]_{i=1, \dots, n, j=1, \dots, r}.$$

Вектор \mathbf{z}_x содержит коэффициенты значимости признаков, \mathbf{z}_y — коэффициенты значимости целевых векторов.

Утверждение 9. Баланс между $\text{Sim}(\mathbf{X})$, $\text{Rel}(\mathbf{X}, \mathbf{Y})$ и $\text{Sim}(\mathbf{Y})$ в задаче (11) достигается при:

$$\alpha_1 \propto \bar{\mathbf{Q}}_y \bar{\mathbf{B}}, \quad \alpha_2 \propto \bar{\mathbf{Q}}_x \bar{\mathbf{Q}}_y, \quad \alpha_3 \propto \bar{\mathbf{Q}}_x \bar{\mathbf{B}}.$$

Здесь $\bar{\mathbf{Q}}_x$, $\bar{\mathbf{B}}$ и $\bar{\mathbf{Q}}_y$ — средние значения соответствующих матриц \mathbf{Q}_x , \mathbf{B} и \mathbf{Q}_y членов $\text{Sim}(\mathbf{X})$, $\text{Rel}(\mathbf{X}, \mathbf{Y})$ и $\text{Sim}(\mathbf{Y})$.

Минимаксная постановка задачи выбора признаков. Сформулируем две взаимосвязанные задачи:

$$\alpha_1 \cdot \underbrace{\mathbf{z}_x^\top \mathbf{Q}_x \mathbf{z}_x}_{\text{Sim}(\mathbf{X})} - \alpha_2 \cdot \underbrace{\mathbf{z}_x^\top \mathbf{B} \mathbf{z}_y}_{\text{Rel}(\mathbf{X}, \mathbf{Y})} \rightarrow \min_{\substack{\mathbf{z}_x \geq \mathbf{0}_n, \\ \mathbf{1}_n^\top \mathbf{z}_x = 1}}, \quad (12)$$

$$\alpha_3 \cdot \underbrace{\mathbf{z}_y^\top \mathbf{Q}_y \mathbf{z}_y}_{\text{Sim}(\mathbf{Y})} + \alpha_2 \cdot \underbrace{\mathbf{z}_x^\top \mathbf{B} \mathbf{z}_y}_{\text{Rel}(\mathbf{X}, \mathbf{Y})} \rightarrow \min_{\substack{\mathbf{z}_y \geq \mathbf{0}_r, \\ \mathbf{1}_r^\top \mathbf{z}_y = 1}}. \quad (13)$$

Задачи (12) и (13) объединяются в совместную минимакс или максмин постановку

$$\min_{\substack{\mathbf{z}_x \geq \mathbf{0}_n \\ \mathbf{1}_n^\top \mathbf{z}_x = 1}} \max_{\substack{\mathbf{z}_y \geq \mathbf{0}_r \\ \mathbf{1}_r^\top \mathbf{z}_y = 1}} f(\mathbf{z}_x, \mathbf{z}_y), \quad \left(\text{или} \max_{\substack{\mathbf{z}_y \geq \mathbf{0}_r \\ \mathbf{1}_r^\top \mathbf{z}_y = 1}} \min_{\substack{\mathbf{z}_x \geq \mathbf{0}_n \\ \mathbf{1}_n^\top \mathbf{z}_x = 1}} f(\mathbf{z}_x, \mathbf{z}_y) \right), \quad (14)$$

где

$$f(\mathbf{z}_x, \mathbf{z}_y) = \alpha_1 \cdot \underbrace{\mathbf{z}_x^\top \mathbf{Q}_x \mathbf{z}_x}_{\text{Sim}(\mathbf{X})} - \alpha_2 \cdot \underbrace{\mathbf{z}_x^\top \mathbf{B} \mathbf{z}_y}_{\text{Rel}(\mathbf{X}, \mathbf{Y})} - \alpha_3 \cdot \underbrace{\mathbf{z}_y^\top \mathbf{Q}_y \mathbf{z}_y}_{\text{Sim}(\mathbf{Y})}.$$

Теорема 2. Для положительно определенной матрицы \mathbf{Q}_x и \mathbf{Q}_y , максмин и минимакс задачи (14) имеют одинаковое оптимальное значение.

Утверждение 10. Минимаксная задача (14) эквивалентна задаче квадратичного программирования с $n + r + 1$ переменными.

Несимметричный учёт значимостей признаков и целевых переменных. Добавим линейный член $\mathbf{b}^\top \mathbf{z}_y$ в член $\text{Rel}(\mathbf{X}, \mathbf{Y})$ следующим образом:

$$\alpha_1 \cdot \underbrace{\mathbf{z}_x^\top \mathbf{Q}_x \mathbf{z}_x}_{\text{Sim}(\mathbf{X})} - \alpha_2 \cdot \underbrace{(\mathbf{z}_x^\top \mathbf{B} \mathbf{z}_y - \mathbf{b}^\top \mathbf{z}_y)}_{\text{Rel}(\mathbf{X}, \mathbf{Y})} + \alpha_3 \cdot \underbrace{\mathbf{z}_y^\top \mathbf{Q}_y \mathbf{z}_y}_{\text{Sim}(\mathbf{Y})} \rightarrow \min_{\substack{\mathbf{z}_x \geq \mathbf{0}_n, \mathbf{1}_n^\top \mathbf{z}_x = 1 \\ \mathbf{z}_y \geq \mathbf{0}_r, \mathbf{1}_r^\top \mathbf{z}_y = 1}}. \quad (15)$$

Утверждение 11. Пусть вектор \mathbf{b} равен

$$b_j = \max_{i=1, \dots, n} [\mathbf{B}]_{i,j}.$$

Тогда значение коэффициентов значимостей вектора \mathbf{z}_y будут неотрицательными в $\text{Rel}(\mathbf{X}, \mathbf{Y})$ для задачи (15).

Утверждение 12. Баланс между членами $\text{Sim}(\mathbf{X})$, $\text{Rel}(\mathbf{X}, \mathbf{Y})$ и $\text{Sim}(\mathbf{Y})$ для задачи (15) достигается при следующих коэффициентах:

$$\alpha_1 \propto \overline{\mathbf{Q}_y} (\overline{\mathbf{b}} - \overline{\mathbf{B}}), \quad \alpha_2 \propto \overline{\mathbf{Q}_x} \overline{\mathbf{Q}_y}, \quad \alpha_3 \propto \overline{\mathbf{Q}_x} \overline{\mathbf{B}}.$$

Теорема 3. В случае скалярной целевой переменной ($r = 1$) предлагаемые методы выбора признаков SymImp (11), MinMax (14), AsymImp (15) совпадают с оригинальным методом QPFS (10).

Таблица 1 демонстрирует основные идеи и функции ошибок для каждого метода. RelAgg является базовой стратегией и не учитывает корреляции в целевом пространстве. SymImp штрафует попарные корреляции между целевыми векторами. MinMax более чувствителен к целевым векторам, которые трудно предсказать. Стратегия Asymimp добавляет линейный член к функции SymImp, чтобы сделать вклад признаков и целевых векторов асимметричным.

В главе 4 предлагается метод выбора активных параметров модели с использованием метода выбора признаков с помощью квадратичного программирования. Приводится анализ параметров модели, которые не находятся в оптимуме.

Таблица 1: Обзор предлагаемых обобщений метода QPFS для векторной целевой переменной

Метод	Идея	Функция ошибки $S(\mathbf{z} \mathbf{X}, \mathbf{Y})$
RelAgg	$\min[\text{Sim}(\mathbf{X}) - \text{Rel}(\mathbf{X}, \mathbf{Y})]$	$\min_{\mathbf{z}_x} [(1 - \alpha) \cdot \mathbf{z}_x^\top \mathbf{Q}_x \mathbf{z}_x - \alpha \cdot \mathbf{z}_x^\top \mathbf{B} \mathbf{1}_r]$
SymImp	$\min [\text{Sim}(\mathbf{X}) - \text{Rel}(\mathbf{X}, \mathbf{Y}) + \text{Sim}(\mathbf{Y})]$	$\min_{\mathbf{z}_x, \mathbf{z}_y} [\alpha_1 \cdot \mathbf{z}_x^\top \mathbf{Q}_x \mathbf{z}_x - \alpha_2 \cdot \mathbf{z}_x^\top \mathbf{B} \mathbf{z}_y + \alpha_3 \cdot \mathbf{z}_y^\top \mathbf{Q}_y \mathbf{z}_y]$
MinMax	$\min [\text{Sim}(\mathbf{X}) - \text{Rel}(\mathbf{X}, \mathbf{Y})]$ $\max [\text{Rel}(\mathbf{X}, \mathbf{Y}) + \text{Sim}(\mathbf{Y})]$	$\min_{\mathbf{z}_x} \max_{\mathbf{z}_y} [\alpha_1 \cdot \mathbf{z}_x^\top \mathbf{Q}_x \mathbf{z}_x - \alpha_2 \cdot \mathbf{z}_x^\top \mathbf{B} \mathbf{z}_y - \alpha_3 \cdot \mathbf{z}_y^\top \mathbf{Q}_y \mathbf{z}_y]$
AsymImp	$\min [\text{Sim}(\mathbf{X}) - \text{Rel}(\mathbf{X}, \mathbf{Y})]$ $\max [\text{Rel}(\mathbf{X}, \mathbf{Y}) + \text{Sim}(\mathbf{Y})]$	$\min_{\mathbf{z}_x, \mathbf{z}_y} [\alpha_1 \cdot \mathbf{z}_x^\top \mathbf{Q}_x \mathbf{z}_x - \alpha_2 \cdot (\mathbf{z}_x^\top \mathbf{B} \mathbf{z}_y - \mathbf{b}^\top \mathbf{z}_y) + \alpha_3 \cdot \mathbf{z}_y^\top \mathbf{Q}_y \mathbf{z}_y]$

Модель $f(\mathbf{x}, \boldsymbol{\theta})$ с параметрами $\boldsymbol{\theta} \in \mathbb{R}^p$ предсказывает целевой объект $y \in \mathbb{Y}$ по исходному объекту $\mathbf{x} \in \mathbb{R}^n$. Пространство \mathbb{Y} представляет собой бинарные метки классов $\{0, 1\}$ для задачи двухклассовой классификации и \mathbb{R} для задачи регрессии. Параметры $\boldsymbol{\theta}$ вычисляются минимизацией функции ошибки:

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^p} \mathcal{L}(\boldsymbol{\theta}, \mathbf{X}, \mathbf{y}). \quad (16)$$

В качестве функции ошибки $\mathcal{L}(\boldsymbol{\theta}, \mathbf{X}, \mathbf{y})$ рассматриваются квадратичная ошибка для задачи регрессии:

$$\mathcal{L}(\boldsymbol{\theta}, \mathbf{X}, \mathbf{y}) = \frac{1}{2} \|\mathbf{y} - \mathbf{f}(\mathbf{X}, \boldsymbol{\theta})\|_2^2 = \frac{1}{2} \sum_{i=1}^m (y_i - f(\mathbf{x}_i, \boldsymbol{\theta}))^2, \quad (17)$$

и функция кросс-энтропии для задачи бинарной классификации:

$$\mathcal{L}(\boldsymbol{\theta}, \mathbf{X}, \mathbf{y}) = \sum_{i=1}^m [y_i \log f(\mathbf{x}_i, \boldsymbol{\theta}) + (1 - y_i) \log(1 - f(\mathbf{x}_i, \boldsymbol{\theta}))]. \quad (18)$$

Для выбора вектора обновлений $\Delta \boldsymbol{\theta}$ используется метод оптимизации Ньютона.

Метод Ньютона нестабилен и вычислительно сложен. В работе предлагается стабильный метод Ньютона. Перед шагом градиента предлагается выбрать подмножество активных параметров модели, которые оказывают наибольшее влияние на функцию ошибки $\mathcal{L}(\boldsymbol{\theta}, \mathbf{X}, \mathbf{y})$.

Определение 11. Параметр θ_j для модели $f(\mathbf{x}, \boldsymbol{\theta})$ является *активным*, если $\mathbf{J}^\top(\mathbf{f}(\mathbf{x}, \boldsymbol{\theta}) - \mathbf{y}) \neq 0$.

Обновление параметров производится только для отобранного множества индексов $\mathcal{A} = \{j : a_j = 1, \mathbf{a} \in \{0, 1\}^p\}$

$$\begin{aligned} \boldsymbol{\theta}_{\mathcal{A}}^k &= \boldsymbol{\theta}_{\mathcal{A}}^{k-1} + \Delta \boldsymbol{\theta}_{\mathcal{A}}^{k-1}, \quad \boldsymbol{\theta}_{\mathcal{A}} = \{\theta_j : j \in \mathcal{A}\}, \\ \boldsymbol{\theta}_{\bar{\mathcal{A}}}^k &= \boldsymbol{\theta}_{\bar{\mathcal{A}}}^{k-1}, \quad \boldsymbol{\theta}_{\bar{\mathcal{A}}} = \{\theta_j : j \notin \mathcal{A}\}. \end{aligned}$$

Чтобы выбрать оптимальное подмножество индексов \mathcal{A} , из всех возможных $2^p - 1$ подмножеств, используется метод QPFS. Метод выбирает подмножество параметров \mathbf{a} для вектора обновлений $\Delta\boldsymbol{\theta}$, которые оказывают наибольшее влияние на вектор остатков и являются попарно независимыми:

$$\mathbf{a} = \arg \max_{\mathbf{a}' \in \{1,0\}^p} S(\mathbf{a}', \mathbf{X}, \mathbf{y}, \boldsymbol{\theta}) \Leftrightarrow \arg \min_{\mathbf{a} \in \mathbb{R}_+^p, \|\mathbf{a}\|_1=1} [\mathbf{a}^\top \mathbf{Q} \mathbf{a} - \alpha \cdot \mathbf{b}^\top \mathbf{a}].$$

Метод Ньютона использует условие оптимизации первого порядка для задачи (16) и линеаризует градиент $S(\boldsymbol{\theta})$

$$\nabla S(\boldsymbol{\theta} + \Delta\boldsymbol{\theta}) = \nabla S(\boldsymbol{\theta}) + \mathbf{H} \cdot \Delta\boldsymbol{\theta} = 0, \quad \Delta\boldsymbol{\theta} = -\mathbf{H}^{-1} \nabla S(\boldsymbol{\theta}).$$

где $\mathbf{H} = \nabla^2 S(\boldsymbol{\theta})$ является матрицей Гессияна функции ошибки $S(\boldsymbol{\theta})$.

Теорема 4. Пусть модель $f(\mathbf{x}, \boldsymbol{\theta})$ близка к линейной в окрестности точки $\boldsymbol{\theta} + \Delta\boldsymbol{\theta}$

$$\mathbf{f}(\mathbf{X}, \boldsymbol{\theta} + \Delta\boldsymbol{\theta}) \approx \mathbf{f}(\mathbf{X}, \boldsymbol{\theta}) + \mathbf{J} \cdot \Delta\boldsymbol{\theta},$$

где $\mathbf{J} \in \mathbb{R}^{m \times p}$ является матрицей Якоби. Тогда вектор обновления $\Delta\boldsymbol{\theta}$ для функции ошибки (17) является решением задачи линейной регрессии

$$\|\mathbf{e} - \mathbf{F} \Delta\boldsymbol{\theta}\|_2^2 \rightarrow \min_{\Delta\boldsymbol{\theta} \in \mathbb{R}^p}, \quad (19)$$

где $\mathbf{e} = \mathbf{f} - \mathbf{y}$ и $\mathbf{F} = \mathbf{J}$.

Теорема 5. Рассмотрим модель логистической регрессии вида $f(\mathbf{x}, \boldsymbol{\theta}) = \sigma(\mathbf{x}^\top \boldsymbol{\theta})$ с сигмоидной функцией активации $\sigma(\cdot)$. Вектор обновлений $\Delta\boldsymbol{\theta}$ для функции ошибки (18) является решением задачи линейной регрессии

$$\|\mathbf{e} - \mathbf{F} \Delta\boldsymbol{\theta}\|_2^2 \rightarrow \min_{\Delta\boldsymbol{\theta} \in \mathbb{R}^p}, \quad (20)$$

где $\mathbf{e} = \mathbf{R}^{-1/2}(\mathbf{y} - \mathbf{f})$ и $\mathbf{F} = \mathbf{R}^{1/2} \mathbf{X}$.

Предлагается адаптация метода QPFS для решения задач (19) и (20). Матрица парных взаимодействий \mathbf{Q} и вектор релевантностей \mathbf{b} имеют вид

$$\mathbf{Q} = \text{Sim}(\mathbf{F}), \quad \mathbf{b} = \text{Rel}(\mathbf{F}, \mathbf{e}).$$

Утверждение 13. В оптимальной точке $\boldsymbol{\theta}^*$ вектор релевантностей $\mathbf{b} = \text{Rel}(\mathbf{F}, \mathbf{e})$ равен нулю.

В главе 5 ставится задача метрического обучения как поиск оптимальной метрики в целевом пространстве. При использовании в качестве функции ошибки квадратичной ошибки предполагается, что целевое пространство является евклидовым. Данное предположение не всегда является адекватным. Рассматриваются задачи кластеризации и классификации множества временных рядов.

Пусть $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_m]^\top \in \mathbb{R}^{m \times n}$ — матрица исходных объектов. Требуется выявить кластерную структуру данных и разбить множество объектов \mathbf{X} на множество непересекающихся кластеров $\mathbb{Y} = \{1, \dots, K\}$, т. е. построить отображение $f: \mathbb{R}^n \rightarrow \mathbb{Y}$. Обозначим $y_i = f(\mathbf{x}_i)$, $y_i \in \mathbb{Y}$ — метка кластера объекта \mathbf{x}_i .

Определение 12. Центроидом класса $e \in \mathbb{Y}$ множества объектов \mathbf{X}_e по расстоянию ρ назовем вектор $\mathbf{c}_e \in \mathbb{R}^n$ такой, что

$$\mathbf{c}_e = \operatorname{argmin}_{\mathbf{c} \in \mathbb{R}^n} \sum_{\mathbf{x}_i \in \mathbf{X}_e} \rho(\mathbf{x}_i, \mathbf{c}). \quad (21)$$

Здесь $\mathbf{X}_e = \{\mathbf{x}_i : i = 1, \dots, m, y_i = e\}$ — множество объектов выборки, принадлежащих одному классу $e \in \mathbb{Y}$.

В случае использования в качестве расстояния ρ евклидовой метрики формула (21) принимает вид

$$\mathbf{c} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i, \quad \mathbf{c}_e = \frac{\sum_{i=1}^m [y_i = e] \mathbf{x}_i}{\sum_{i=1}^m [y_i = e]}.$$

Здесь вектор \mathbf{c} является центроидом всего множества объектов \mathbf{X} .

Введем на множестве объектов \mathbf{X} расстояние Махаланобиса

$$d_{\mathbf{A}}(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{(\mathbf{x}_i - \mathbf{x}_j)^\top \mathbf{A}^{-1} (\mathbf{x}_i - \mathbf{x}_j)}, \quad (22)$$

где матрица трансформаций $\mathbf{A} \in \mathbb{R}^{n \times n}$ является симметричной и неотрицательно определенной ($\mathbf{A}^\top = \mathbf{A}$, $\mathbf{A} \succeq 0$). Зададим в качестве матрицы трансформации матрицу выборочной ковариации

$$\mathbf{A} = \frac{1}{m} \sum_{i=1}^m (\mathbf{x}_i - \mathbf{c})(\mathbf{x}_i - \mathbf{c})^\top. \quad (23)$$

Функцией ошибки кластеризации назовем межкластерное расстояние:

$$\mathcal{L}(\{\mathbf{c}_j\}_{j=1}^K, \mathbf{X}, \mathbf{y}) = - \sum_{j=1}^K N_j d_{\mathbf{A}}^2(\mathbf{c}_j, \mathbf{c}), \quad (24)$$

где $N_j = \sum_{i=1}^m [y_i = y_j]$ — число объектов в кластере j .

Поставим задачу кластеризации как задачу минимизации функции ошибки (24)

$$\mathcal{L}(\{\mathbf{c}_j\}_{j=1}^K, \mathbf{X}, \mathbf{y}) \rightarrow \min_{\mathbf{c}_j \in \mathbb{R}^n}. \quad (25)$$

Найдем такую матрицу \mathbf{A} , для которой функционал качества принимает максимальное значение:

$$\mathbf{A}^* = \arg \min_{\mathbf{A} \in \mathbb{R}^{n \times n}} S(\{\mathbf{c}_j^*\}_{j=1}^K, \mathbf{X}, \mathbf{y}), \quad (26)$$

где $\{\mathbf{c}_j^*\}_{j=1}^K$ — решение задачи кластеризации (25).

Для решения задач (25), (26) используется алгоритм адаптивного метрического обучения. Предлагается понизить размерность пространства объектов \mathbf{X}

с помощью линейного ортогонального преобразования $\mathbf{P} \in \mathbb{R}^{l \times n}$, $\mathbf{P}^\top \mathbf{P} = \mathbf{I}$, где новая размерность $l < n$

$$\mathbf{t}_i = \mathbf{P} \mathbf{x}_i \in \mathbb{R}^l, \quad i = 1, \dots, m.$$

Расстояния между объектами вычисляются по формуле (22), где в качестве матрицы $\hat{\mathbf{A}}$ используется матрица ковариаций (23) множества объектов $\{\hat{\mathbf{x}}_i\}_{i=1}^m$

$$\hat{\mathbf{A}} = \frac{1}{m} \sum_{i=1}^m (\mathbf{t}_i - \hat{\mathbf{c}})(\mathbf{t}_i - \hat{\mathbf{c}})^\top = \frac{1}{m} \sum_{i=1}^m \mathbf{P}(\mathbf{x}_i - \mathbf{c})(\mathbf{x}_i - \mathbf{c})^\top \mathbf{P}^\top = \mathbf{P} \mathbf{A} \mathbf{P}^\top.$$

Определение 13. *Индикаторной матрицей* назовем матрицу $\mathbf{Y} \in \mathbb{R}^{m \times K}$, где $y_{ij} = [f(\mathbf{x}_i) = y_j]$.

Определение 14. *Взвешенной индикаторной матрицей* назовем матрицу $\mathbf{L} = \mathbf{Y}(\mathbf{Y}^\top \mathbf{Y})^{-1/2} \in \mathbb{R}^{m \times K}$, элементы которой равны:

$$l_{ij} = \begin{cases} \frac{1}{\sqrt{N_j}}, & \text{если } f(\mathbf{x}_i) = y_j; \\ 0, & \text{если } f(\mathbf{x}_i) \neq y_j. \end{cases}$$

Задача кластеризации (25) и задача метрического обучения (26) сводятся к общей задаче минимизации функции ошибки

$$\begin{aligned} \mathcal{L} &= -\frac{1}{m} \text{trace}(\mathbf{L}^\top \mathbf{X}^\top \mathbf{P}^\top \hat{\mathbf{A}}^{-1} \mathbf{P} \mathbf{X} \mathbf{L}) = \\ &= -\frac{1}{m} \text{trace}(\mathbf{L}^\top \mathbf{X}^\top \mathbf{P}^\top (\mathbf{P} \mathbf{A} \mathbf{P}^\top)^{-1} \mathbf{P} \mathbf{X} \mathbf{L}) \rightarrow \min_{\mathbf{P}, \mathbf{L}}. \end{aligned} \quad (27)$$

Для решения задачи (27) используется ЕМ алгоритм.

Оптимизация матрицы \mathbf{P} с фиксированной матрицей \mathbf{L} . Переформулируем задачу (27) следующим образом:

$$\mathcal{L} = -\frac{1}{m} \text{trace}(\mathbf{L}^\top \mathbf{X}^\top \mathbf{P}^\top (\mathbf{P} \mathbf{A} \mathbf{P}^\top)^{-1} \mathbf{P} \mathbf{X} \mathbf{L}) = -\frac{1}{m} \text{trace}((\mathbf{P} \mathbf{A} \mathbf{P}^\top)^{-1} \mathbf{P} \mathbf{X} \mathbf{L} \mathbf{L}^\top \mathbf{X}^\top \mathbf{P}^\top).$$

Утверждение 14. Обозначим $\mathbf{B} = \mathbf{X} \mathbf{L} \mathbf{L}^\top \mathbf{X}^\top$. Обозначим через $\mathbf{P} = [\mathbf{v}_1, \dots, \mathbf{v}_r]^\top$ матрицу, состоящую из r собственных векторов матрицы $\mathbf{A}^{-1} \mathbf{B}$, отвечающих наибольшему собственному значению. Тогда решением (27) является ортогональная матрица, полученная QR-разложением матрицы \mathbf{P}^\top .

Оптимизация матрицы \mathbf{L} с фиксированной матрицей \mathbf{P} . Обозначим $\hat{\mathbf{K}} = (1/N) \mathbf{X}^\top \mathbf{P}^\top \hat{\mathbf{A}}^{-1} \mathbf{P} \mathbf{X}$.

При фиксированной матрице \mathbf{P} задача (27) принимает вид:

$$\text{trace}(\mathbf{L}^\top \hat{\mathbf{K}} \mathbf{L}) \rightarrow \max_{\mathbf{L} \in \mathbb{R}^{m \times r}}.$$

Для нахождения оптимального соответствия между временными рядами при решении задачи классификации предлагается процедура выравнивания временных рядов. Пусть исходный объект $\mathbf{x}_i \in \mathbb{R}^n$ — временной ряд, последовательность измерений некоторой исследуемой величины в различные моменты времени. Пусть задана выборка $\{(\mathbf{x}_i, y_i)\}_{i=1}^m$ — множество объектов с известными метками классов $y_i \in \mathbb{Y}$, где $\mathbb{Y} = \{1, \dots, K\}$ — множество меток классов.

Требуется построить точную, простую, устойчивую модель классификации $a : \mathbb{R}^n \rightarrow \mathbb{Y}$. Данную модель представим в виде суперпозиции

$$a(\mathbf{x}) = b \circ \mathbf{f} \circ G(\mathbf{x}, \{\mathbf{c}_e\}_{e=1}^K),$$

где G — процедура выравнивания временных рядов относительно центроидов классов $\{\mathbf{c}_e\}_{e=1}^K$, \mathbf{f} — алгоритм метрического обучения, b — алгоритм многоклассовой классификации.

Для нахождения центроида предлагается в качестве расстояния ρ между временными рядами использовать путь наименьшей стоимости $[?, ?]$, найденный методом динамической трансформации времени.

Метрическое обучение. Введем на множестве выравненных временных рядов расстояние Махаланобиса $d_{\mathbf{A}}$ 22. Представим матрицу трансформации \mathbf{A} в виде разложения $\mathbf{A}^{-1} = \mathbf{L}^T \mathbf{L}$. Матрица $\mathbf{L} \in \mathbb{R}^{p \times n}$ — матрица линейного преобразования, где p задает размерность преобразованного пространства.

Расстояние $d_{\mathbf{A}}(\mathbf{x}_i, \mathbf{x}_j)$ есть евклидово расстояние между $\mathbf{L}\mathbf{x}_i$ и $\mathbf{L}\mathbf{x}_j$:

$$\begin{aligned} d_{\mathbf{A}}(\mathbf{x}_i, \mathbf{x}_j) &= \sqrt{(\mathbf{x}_i - \mathbf{x}_j)^T \mathbf{L}^T \mathbf{L} (\mathbf{x}_i - \mathbf{x}_j)} = \\ &= \sqrt{(\mathbf{L}(\mathbf{x}_i - \mathbf{x}_j))^T (\mathbf{L}(\mathbf{x}_i - \mathbf{x}_j))} = \|\mathbf{L}(\mathbf{x}_i - \mathbf{x}_j)\|_2. \end{aligned}$$

В качестве алгоритма метрического обучения в данной работе был выбран алгоритм LMNN. Данный алгоритм сочетает в себе идеи метода k ближайших соседей. Первая идея заключается в минимизации расстояний между k ближайшими объектами, находящимися в одном классе:

$$Q_1(\mathbf{L}) = \sum_{j \rightsquigarrow i} \|\mathbf{L}(\mathbf{x}_i - \mathbf{x}_j)\|^2 \rightarrow \min_{\mathbf{L}},$$

где $j \rightsquigarrow i$ означает, что \mathbf{x}_j является одним из k ближайших соседей для \mathbf{x}_i . Вторая идея состоит в максимизации расстояния между каждым объектом и его объектами-нарушителями.

Определение 15. *Объектом-нарушителем* для \mathbf{x}_i назовем объект \mathbf{x}_l такой, что

$$\|\mathbf{L}(\mathbf{x}_i - \mathbf{x}_l)\|^2 \leq \|\mathbf{L}(\mathbf{x}_i - \mathbf{x}_j)\|^2 + 1, \quad \text{где } j \rightsquigarrow i.$$

Таким образом, необходимо минимизировать следующий функционал:

$$Q_2(\mathbf{L}) = \sum_{j \rightsquigarrow i} \sum_l [y_i \neq y_l] [1 + \|\mathbf{L}(\mathbf{x}_i - \mathbf{x}_j)\|^2 - \|\mathbf{L}(\mathbf{x}_i - \mathbf{x}_l)\|^2]_+ \rightarrow \min_{\mathbf{L}}.$$

Задача метрического обучения состоит в нахождении линейного преобразования $\mathbf{f}(\mathbf{x}) = \mathbf{L}\mathbf{x}$, то есть нахождении матрицы \mathbf{L} в виде решения оптимизационной задачи

$$Q(\mathbf{L}) = \mu Q_1(\mathbf{L}) + (1 - \mu) Q_2(\mathbf{L}) \rightarrow \min_{\mathbf{L}}.$$

В **главе 6** ставится задача порождения информативного признакового пространства. Временные ряды акселерометра образуют множество \mathcal{S} сегментов $\mathbf{s} = [x_1, \dots, x_T]^\top$ фиксированной длины T . Необходимо построить модель классификации $a : \mathbb{R}^T \rightarrow \mathbb{Y}$, которая будет ставить в соответствие каждому сегменту из множества \mathcal{S} метку класса из конечного множества $\mathbb{Y} = \{1, \dots, K\}$. Обозначим за $\mathcal{D} = \{(\mathbf{s}_i, y_i)\}_{i=1}^m$ исходную выборку, где $\mathbf{s}_i \in \mathcal{S}$ и $y_i = a(\mathbf{s}_i) \in \mathbb{Y}$.

В работе предлагается построить модель a в виде суперпозиции $a = f \circ \mathbf{g}$. **Определение 16.** *Порождающей функцией* будем называть функцию $\mathbf{g} : \mathbb{R}^T \rightarrow \mathbb{X}$, отображающую исходные временные ряды \mathbf{s} из пространства \mathbb{R}^T в признаковое пространство $\mathbb{X} \subset \mathbb{R}^n$.

Имея порождающую функцию \mathbf{g} , преобразуем исходную выборку в $\mathcal{D}_{\mathbb{X}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$, где $\mathbf{x}_i = \mathbf{g}(\mathbf{s}_i) \in \mathbb{X}$. Здесь \mathbf{x}_i является исходным объектом.

Модель классификации $f = f(\mathbf{x}, \boldsymbol{\theta})$ является параметрической функцией с вектором параметров $\boldsymbol{\theta}$. Оптимальные параметры $\boldsymbol{\theta}^*$ определяются оптимизацией функции ошибки классификации

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} L(\boldsymbol{\theta}, \mathcal{D}_{\mathbb{X}}, \mu).$$

Экспертные функции. Экспертные функции — это некоторые статистики g_j , где $g_j : \mathbb{R}^T \rightarrow \mathbb{R}$. Признаковым описанием $\mathbf{g}(\mathbf{s})$ объекта \mathbf{s} являются значения заданных экспертных статистик для данного объекта

$$\mathbf{x} = \mathbf{g}(\mathbf{s}) = [g_1(\mathbf{s}), \dots, g_n(\mathbf{s})]^\top.$$

Авторегрессионная модель. Авторегрессионная модель порядка n использует параметрическую модель для аппроксимации временного ряда \mathbf{s} :

$$x_t = w_0 + \sum_{j=1}^{n-1} w_j x_{t-j} + \varepsilon_t.$$

Оптимальные параметры \mathbf{w}^* авторегрессионной модели используются как признаки $\mathbf{g}(\mathbf{s})$:

$$\mathbf{x} = \mathbf{g}(\mathbf{s}) = \mathbf{w}^* = \arg \min_{\mathbf{w} \in \mathbb{R}^n} \left(\sum_{t=n}^T \|x_t - \hat{x}_t\|^2 \right).$$

Таблица 2: Примеры экспертных порождающих функций

Описание	Формула
Mean	$\bar{s} = \frac{1}{T} \sum_{t=1}^T s_t$
Standard deviation	$\sqrt{\frac{1}{T} \sum_{t=1}^T (s_t - \bar{s})^2}$
Mean absolute deviation	$\frac{1}{T} \sum_{t=1}^T s_t - \bar{s} $
Distribution	Histogram values with 10 bins

Анализ сингулярного спектра. Для каждого временного ряда \mathbf{s} из исходной выборки \mathcal{D} строится траекторная матрица:

$$\mathbf{S} = \begin{pmatrix} s_1 & s_2 & \dots & s_n \\ s_2 & s_3 & \dots & s_{n+1} \\ \dots & \dots & \dots & \dots \\ s_{T-n+1} & s_{T-n+2} & \dots & s_T \end{pmatrix}.$$

Сингулярное разложение матрицы $\mathbf{S}^\top \mathbf{S}$:

$$\mathbf{S}^\top \mathbf{S} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^\top.$$

Признаковое описание объекта \mathbf{s} задаётся спектром матрицы $\mathbf{S}^\top \mathbf{S}$:

$$\mathbf{x} = \mathbf{g}(\mathbf{s}) = [\lambda_1, \dots, \lambda_n]^\top.$$

Аппроксимация сплайнами. Предполагается, что узлы сплайна $\{\xi_\ell\}_{\ell=0}^M$ равномерно распределены по временной оси. Кусочные модели, построенные на отрезках $[\xi_{\ell-1}; \xi_\ell]$, заданы коэффициентами $\{\mathbf{w}_\ell\}_{\ell=1}^M$. Обозначим каждый отрезок-сегмент $p_i(t)$ $i = 1, \dots, M$ и весь сплайн $S(t)$:

$$S(t) = \begin{cases} p_1(t) = w_{10} + w_{11}t + w_{12}t^2 + w_{13}t^3, & t \in [\xi_0, \xi_1], \\ p_2(t) = w_{20} + w_{21}t + w_{22}t^2 + w_{23}t^3, & t \in [\xi_1, \xi_2], \\ \dots & \dots \\ p_M(t) = w_{L0} + w_{M1}t + w_{M2}t^2 + w_{M3}t^3, & t \in [\xi_{M-1}, \xi_M], \end{cases}$$

$$\begin{aligned} S(\xi_t) &= x_t, \quad t = 0, \dots, M, \\ p'_i(\xi_i) &= p'_{i+1}(\xi_i), p''_i(\xi_i) = p''_{i+1}(\xi_i), \quad i = 1, \dots, M-1, \\ p_i(\xi_{i-1}) &= x_{i-1}, p_i(\xi_i) = x_i, \quad i = 1, \dots, M. \end{aligned}$$

Объединение всех параметров сплайна задаёт признаковое описание временного ряда:

$$\mathbf{x} = \mathbf{g}(\mathbf{s}) = [\mathbf{w}_1, \dots, \mathbf{w}_M]^\top.$$

В заключении представлены основные результаты диссертационной работы.

1. Исследована проблема снижения размерности сигналов в коррелированных пространствах высокой размерности. Предложены методы декодирования сигналов, учитывающие зависимости как в исходном, так и в целевом пространстве сигналов.
2. Доказаны теоремы об оптимальности предлагаемых методов декодирования сигналов. Предлагаемые методы выбирают согласованные модели в случае избыточной размерности описания данных.
3. Предложены методы выбора признаков, учитывающие зависимости как в исходном, так и в целевом пространстве. Предложенные методы доставляют устойчивые и адекватные решения в пространствах высокой размерности.
4. Предложены нелинейные методы согласования скрытых пространств для данных со сложноорганизованной целевой переменной.
5. Предложен ряд моделей для прогнозирования гетерогенных наборов сигналов для задачи построения нейрокомпьютерных интерфейсов.

Публикации соискателя по теме диссертации

Публикации в журналах из списка ВАК.

1. Исаченко Р. В., Стрижов В. В. Метрическое обучение в задачах мультиклассовой классификации временных рядов // Информатика и её применения, 2016. Т. 10. № 2. С. 48–57.
2. Isachenko R. et al. Feature Generation for Physical Activity Classification // Artificial Intelligence and Decision Making, 2018. № 3. С. 20–27.
3. Isachenko R. V., Strijov V. V. Quadratic programming optimization with feature selection for nonlinear models // Lobachevskii Journal of Mathematics, 2018. Т. 39. № 9. С. 1179–1187.
4. Isachenko R. V., Vladimirova M. R., Strijov V. V. Dimensionality Reduction for Time Series Decoding and Forecasting Problems // DEStech Transactions on Computer Science and Engineering, 2018. №. optim.
5. Исаченко Р.В., Яушев Ф.Ю., Стрижов В.В. Модели согласования скрытого пространства в задаче прогнозирования // Системы и средства информатики, 2021. Т. 31. № 1.

Прочие публикации.

6. Исаченко Р. В., Катруца А. М. Метрическое обучение и снижение размерности пространства в задачах кластеризации // Машинное обучение и анализ данных, 2016. Т. 2. № 1. С. 17–25.